

Scaling behavior of online human activity

ZHI-DAN ZHAO¹, SHI-MIN CAI^{1(a)}, JUNMING HUANG², YAN FU¹, and TAO ZHOU¹

¹ *Web Sciences Center, School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu 610054, P. R. China*

² *Institute of Computing Technology, Chinese Academy of Sciences, Beijing, 100190, P. R. China*

PACS 89.75.Da – Systems obeying scaling laws

PACS 05.45.Tp – Time series analysis

PACS 89.65.-s – Social and economic systems

Abstract –The rapid development of Internet technology enables human explore the web and record the traces of online activities. From the analysis of these large-scale data sets (i.e. traces), we can get insights about dynamic behavior of human activity. In this letter, the scaling behavior and complexity of human activity in the e-commerce, such as music, book, and movie rating, are comprehensively investigated by using detrended fluctuation analysis technique and multiscale entropy method. Firstly, the interevent time series of rating behaviors of these three type medias show the similar scaling property with exponents ranging from 0.53 to 0.58, which implies that the collective behaviors of rating media follow a process embodying self-similarity and long-range correlation. Meanwhile, by dividing the users into three groups based their activities (i.e., rating per unit time), we find that the scaling exponents of interevent time series in three groups are different. Hence, these results suggest the stronger long-range correlations exist in these collective behaviors. Furthermore, their information complexities vary from three groups. To explain the differences of the collective behaviors restricted to three groups, we study the dynamic behavior of human activity at individual level, and find that the dynamic behaviors of a few users have extremely small scaling exponents associating with long-range anticorrelations. By comparing with the interevent time distributions of four representative users, we can find that the bimodal distributions may bring the extraordinary scaling behaviors. These results of analyzing the on-line human activity in the e-commerce may not only provide insights to understand its dynamic behaviors but also be applied to acquire the potential economic interest.

Introduction. – The human behavior involving their daily activities is one of the highest complexity and complicated things because it is driven by countless unknown facts. Mining the human dynamics from these recorded large-scale data sets has become much more important for understanding their behavior patterns, and modeling the human dynamic behaviors helps us to explain many socioeconomic phenomena and find significant applications ranging from resource allocation, transportation control, epidemic prediction to personalized recommendation [1,2]. Thanks to the development of information technology, massive Internet data and resources make us easily realize the empirical analysis and modeling of human activity. One of the most attractive observation is the heavy-tailed nature of the interevent time distribution, which implies that the bursts of rapidly occurring events are separated

by long periods of inactivity. Examples of empirical studies include the email communication [3], the surface mail communication [4], the cell-phone communication [5], the online activities [6–10], and so on. To understand the heavy-tailed phenomena of human dynamic behavior, the experts have put forward many mechanisms, such as the highest-priority-first queue model [3], the varying interest [11], the memory effects [12] and the human interacting [13–15] to mimic the temporal bursts.

On the other hand, the techniques of time series analysis are applied to investigate the evolutionary data in real world. One of the most popular techniques is detrended fluctuation analysis (DFA) proposed by Peng *et al* [16,17], which can effectively quantify the long-range power-law correlations embedded in the nonstationary time series (or self-similarity process). It provides a simple scaling exponent α to represent the correlation characters of time se-

^(a)shimin.cai81@gmail.com

ries, and thus is applied to various research fields including heart rate dynamics [17–21], financial time series [22–26], particle condensation [27], Internet traffic [28], musical rhythm spectra [29], etc. In recent, Rybski *et al* [30–32] have applied the DFA to study the long-term correlations of the communication patterns (i.e., interactive activities) in online social networks, which is associated with the source of general Gibrat’s law in economics. Meanwhile, Costa *et al.* [33] recently proposed a technique, namely multiscale entropy (MSE), to quantify the information complexity of physiologic time series over multiple scale. They further used the MSE to analyze the human heartbeat [34], which suggested that the time asymmetry is a fundamental property of healthy. In the following works, the MSE was widely applied to analyze the EEG signals, which indicated that the human brain variability increases with maturation [35] and the Alzheimer disease patients usually had lower sample entropy on the small and medium time scales [36]. In the environmental field, Li and Zhang [37] analyzed the long-term daily flow rates of the Mississippi River, and found that the sample entropy for flow rates generally monotonously increases with scale factor and the complexity was beginning to decrease since 1940s. Moreover, based on the interevent time distributions and memory, Goh *et al* [38] used the orthogonal measures to quantify the burstiness in many real interevent time series, and found that the origin of burstiness in human activity was much more correlated with the changes in the interevent time distributions.

The e-commerce is composed of the online business trades among humans and rating information based on the Internet Technology, in which the dynamic behaviors of human activity (i.e, trading or rating records) involve with a large number of useful knowledge for acquiring potential economic interest. In this letter, we first focus on the interevent intervals (i.e., time series) of human activity in e-commerce, and empirically investigate their scaling properties and information complexities both in the collective and individual levels based on the measures including the DFA, MSE and interevent time distribution. The rich results include that (1) The interevent time series of rating behaviors restricted to the types of medias show the similar scaling property that implies the collective behaviors of rating media follow a process embodying self-similarity and long-range correlation. (2) The different scaling exponents of interevent time series can be observed from the collective behaviors restricted to users’ activities (i.e., rating per unit time), yet they both suggest the stronger long-range correlations existing in these collective behaviors. (3) The information complexities of collective behaviors are obviously distinguishable with the users’ activities. (4) the extremely small scaling exponents (indicating long-range anticorrelation) of a few representative individual users are mainly brought by the bimodal interevent time distributions.

Materials. – The experimental data set is randomly sampled from Douban, which is a companion of e-commerce. It is similar to the Social Networking Services (SNS) that allows registered users to record information and create content related to movies, books, and music, yet it also can make a personalized recommendation for the registered users. We focus on users who perform more than 1000 rating actions on all three types of medias, which results in a set of 65 individuals. In the data set, we can find series of important history knowledge of registered users, such as user ID, item ID, rate, time stamp and item type, etc. Note that the sample time resolution is second, and we here focus on the interevent time series defined as the intervals between two consecutive rating actions.

Methods. – Herein we apply the DFA and MSE methods to quantitatively understand the scaling behavior and complexity of human activity in e-commerce. In order to keep our description as self-contained as possible, we should introduce the DFA and MSE methods briefly.

Detrended Fluctuation Analysis. – We describe the process of DFA which involves the following steps [16, 39, 40].

(i) Starting with a time series $t(i)$, where $i \in [1, N]$ and N is the length of the series. We first integrate the series $t(i)$ and obtain $y(k) \equiv \sum_{i=1}^k [t(i) - \langle t \rangle]$, where $\langle t \rangle$ is the mean. Meanwhile, $y(k)$ is divided into N/l nonoverlapping boxes, each containing l interevent intervals.

(ii) In the k -th box, we use a polynomial function $y_l(k)$ of order n to represents the local trend. In the experiments, the order is selected as $n = 2$, and the algorithm is denoted as DFA-2.

(iii) We calculate the variance of residual time series after the detrending procedure,

$$F(l) \equiv \sqrt{\frac{1}{N} \sum_{k=1}^N [y(k) - y_l(k)]^2} \quad (1)$$

(iv) Altering the box size l and repeating the detrending procedure, we can obtain the variances $F(l)$ as a function of box size l . A power-law relations between $F(l)$ and l is $F(l) \sim l^\alpha$. The α is a real value in the bounded range from 0 to 1, where $\alpha > 0.5$ means that the time series is correlated, $\alpha = 0.5$ suggests that the time series is same to the white noise (i.e., no correlation), and $\alpha < 0.5$ indicates that the time series is anticorrelated.

Multiscale Entropy. – We look back over the MSE method. The MSE is based on the simple observation that the complex signals generally exhibit the dynamics deviating far from perfect regularity and their multiscale complexity. The procedure of MSE is described as follows [33]:

(i) For a given time series, x_1, x_2, \dots, x_N , where N is the total number of time series. we divide it into nonoverlapping boxes with the length l .

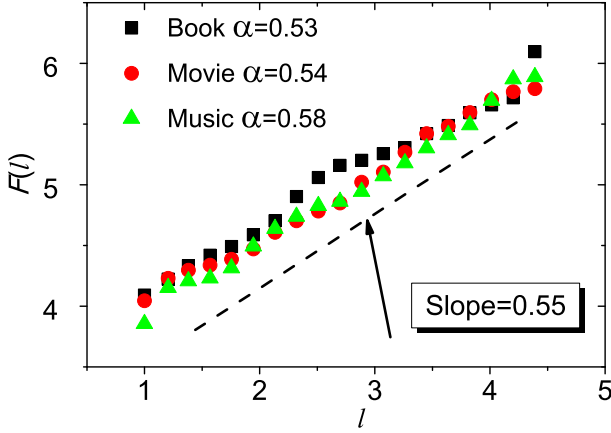


Fig. 1: (Color online) The results of interevent time series measured by DFA. The symbols represent the types of medias, book (black squares), movie (red circles), and music (green triangles). The dash line is presented as guide eyes line. The similar scaling behaviors suggest that the human activity in system follow a in a long-correlated self-similar process.

(ii) The averages of time series inside each boxes are deemed as the elements of a coarse-grained time series,

$$s_j^l = \frac{1}{l} \sum_{i=(j-1)l+1}^{jl} x_i, \quad 1 \leq j \leq N/l. \quad (2)$$

By altering the box length l , we can obtain many coarse-grained time series, which characterize the original time series at multiple scales.

(iii) The sample entropy [41] is used to measure each coarse-grained time series. Thus, we can find the relation between the entropy measure and scale factor (i.e., the box length l).

Collective Level. — The timestamp of data set is in precision of one second. Our focus is the interevent interval τ between consecutive actions, i.e. rating a media by a certain user in Douban. The interevent time series are composed of these intervals in many ways. From the view of whole system, we first investigate the interevent time series restricted to the types of medias. Figure 1 shows that the results of interevent time series measured by DFA, in which we observe that the scaling exponent α fluctuates in the small interval [0.53, 0.58]. The values of α reveal similar scaling behaviors for all types of medias and suggest that the interevent time series of rating media in the system evolves a process embodying self-similarity and long-range correlation. Furthermore, the scaling exponents slightly more than 0.5 also imply the weak memory of the signals, which is consistent with the previous results found in other human activity [38].

Although the human activity in whole system approximately obeys a common scaling law, we should pay attention on the effect of individual user activity on their

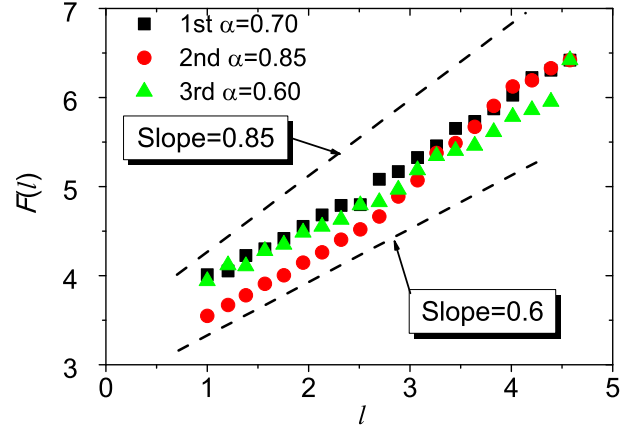


Fig. 2: (Color online) The results of interevent time series measured by DFA. The interevent time series are constructed from the intervals between consecutive actions in the three groups, 1st (black squares), 2nd (red circles), and 3rd (green triangles), respectively. The scaling exponents are 0.7, 0.85, and 0.65 corresponding to 1st, 2nd, and 3rd groups respectively.

scaling behaviors because the user activity is strongly associated with the well understanding of human dynamics [7, 10]. The activity of an arbitrary user i , is defined as $A_i = n_i/T_i$, where n_i is the number of actions and T_i is the time difference between the first and last actions [42]. We sort these users in an increasing order according to their activities, and then divide them into three groups which are indicated by low (1st), mid (2nd) and high (3rd) activity. The number of users in each group is 21 (1st), 22 (2nd), and 22 (3rd). The interevent time series are constructed from the intervals between consecutive actions done by users in the three groups, respectively. Their results measured by DFA are shown in Fig. 2, which suggests that the scaling behaviors are different from the groups. For the 2nd group, the scaling exponent $\alpha = 0.85$ indicates that the interevent time series opposes much stronger long-range correlations than these of other two groups. However, all the scaling exponents are much larger than 0.5, which suggests that the long-range correlations generally exist in the interevent time series regardless of their activities.

Additionally, to verify the the memory effects of online human activities don't arise from the pow law distributions of interevent time series, we reshuffle the original interevent time series to disturb their long-range correlations, and measure these shuffled ones by DFA method. In Tab. 1, we show the remarkable differences of scaling exponents between the original and shuffled time series, which suggest that the memory effects generally exist in online human activities.

To further determine the differences among the three groups, we use the MSE method to quantify the information complexities of interevent time series. There are

Table 1: The scaling exponent α of original and shuffled time series measured by DFA. Specifically, the first column corresponds to the scaling exponents α of original time series, and the second column indicates those of shuffled time series, respectively. Their difference is remarkable, which verifies that the memory effects generally exist in the online human activities.

	<i>Original α</i>	<i>Shuffled α</i>
Book	0.53	0.51
Movie	0.54	0.50
Music	0.58	0.50
1st	0.70	0.50
2nd	0.85	0.52
3rd	0.60	0.51

two guidelines, the higher information complexity is in correspondence with the larger sample entropy and the monotonic increase of the sample entropy indicates the much more information of interevent time series at large scale factors, to compare them. Figure 3 shows that the sample entropy increases with the growth of user activity at all scale factors, and for all interevent time series they first increase at small scale factors and then approximately stabilize at the constant values, which suggests that the interevent time series constructed from the more active users become more complex and contain the much more information at large scale factors. These results demonstrate that the more actions done by users lead to the relative homogeneous interevent time series (i.e., the less extreme intervals). We should notice that although the information complexities are different from the three groups, yet they don't directly associate with the degrees of long-range correlation.

Individual Level. — Although the obvious differences exist in the users behaviors at collective level, we still need to understand them and explore the underlying mechanism of individual user behavior. Therefore, the interevent time series constructed from the total 65 individual user behaviors are further investigated by DFA. The users belonging to different groups are denoted by the different symbols. In Fig. 4, the scaling exponents are as a function of the activity of user, which shows that there is no direct correlation between them. Furthermore, most of scaling exponents greater than 0.5 once demonstrate the existence of long-range correlations in the interevent time series of human activity in e-commerce. However, there are also abnormal user behaviors (e.g., users A and B seen in Fig. 4) suggested by the scaling exponents much less than 0.5.

This phenomenon urges us to observe and study these interevent time series constructed from the abnormal user behavior. Table 2 shows the basic statistical features of abnormal users A and B as well as these of two normal users C and D. From the Tab. 2, we can observe that the the similar activities may show the completely differ-

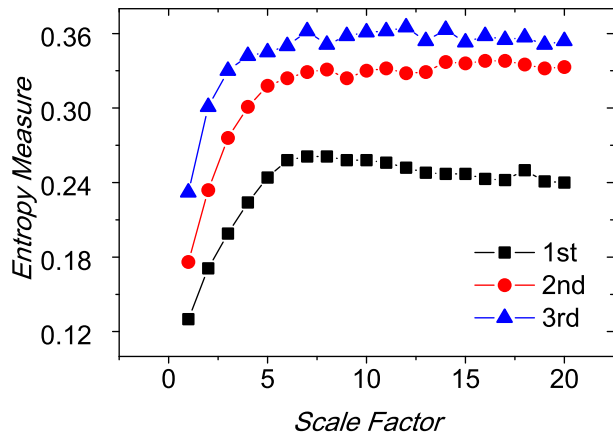


Fig. 3: (Color online) The results of interevent time series measured by MSE. The sample entropy are obviously different from the activity of user behaviors, and much more information is contained at large scale factors. The symbols indicate 1st (black squares), 2nd (red circles) and 3rd (blue triangles) groups, respectively.

ent scaling behaviors (e.g., users B and C) and the similar scaling behaviors don't mean that they have same the activities (e.g., users C and D). We note that the Frequency denotes the event number of user behavior in Tab.1.

To uncover the origin of the observed differences among scaling behaviors of individual users, we first present the rating activity evolving with time in Fig. 5. We can straightly find that the active patterns of users A and B are very different from users C and D. Concretely, the much more actions are occurred in users A and B at the initial stage, which results in the shorter time intervals, and the actions become much less or even absent (e.g., user A) when the time evolves.

We propose the query on what is the specific active pattern for four users, and therefore statistically illustrates the interevent time distributions of the four users in Fig. 6, where the interevent time series are defined as the time intervals between consecutive actions by a certain user. The results confirm that the interevent time distributions of these online human activities follow a power law form, $p(\tau) \sim \tau^{-\beta}$. However, we should note that there is a cutoff at the minute scale for the the interevent time distribution of users A and B, respectively, while for those of users C and D, the blurry cutoffs are at the day scale. The scales of minute and day are a typical decay length of online interests, for example, user actions usually appear within a day in the microblogging systems [10]. Based on the phenomena of bimodal interevent time distributions, it can be found that the power-law exponents β of users A and B change from big to small, yet these cases are quite reverse for users C and D. These changing trends of the bimodal interevent time distributions vividly describe the differences among active patterns of individual user behavior, which suggests that the short and long in-

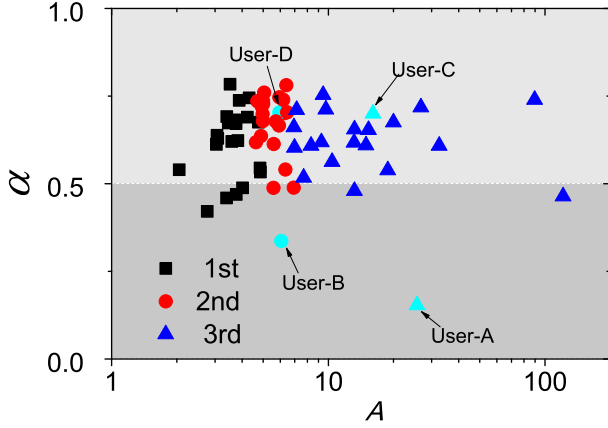


Fig. 4: (Color online) A scatter plot shows the detrended fluctuation analysis of the time series of different users' activities. Each point corresponds to a different user, indicating that there are significant differences between the scaling exponents and users' activity. Where "black squares", "red circles" and "blue triangles" denote the individuals in 1st, 2nd and 3rd activity group, respectively.

terevent intervals for users A and B alternately occur and the short (or long) interevent intervals for users C and D continuously emerge. Thus, we think that the scaling behaviors strongly associate with these changing trends of interevent time distributions for human activity in e-commerce although the potential dynamic mechanisms of online individual activity are similar.

Conclusions. — We conclude that our empirically analysis, including the scaling behaviors and information complexities of human activity (i.e., rating the medias including music, book, and movie) comprehensively investigated by using DFA and MSE methods, provides the well understanding of behavior patterns of human activity in e-commerce. We also find that, for all rating behaviors corresponding to the types of medias, they display the similar scaling property with exponents ranging from 0.53 to 0.58, which implies that the collective behavior pattern of rating media follows a process embodying self-similarity and long-range correlation. Furthermore, by dividing the users into three groups based on their activity, we observe that the scaling exponents among three groups are different, yet they both suggest the stronger long-range correlations exist in the collective behaviors. Meanwhile, the information complexities of human activity quantified by MSE confirm the differences of scaling behaviors in these three groups. Moreover, we study the behavior patterns of human activity at individual level, and find that the behaviors of a few users have extremely small scaling exponents associating with long-range anticorrelations. By comparing with the distributions of interevent time of four representative users, we think that the different scaling behaviors are brought by the bimodal forms of the interevent time distributions.

Table 2: The basic statistical features of the four selected typical users. The first column corresponds to the scaling exponents α , the second column indicates the activities, and the third column represents the frequency, respectively.

	α	Activity(day)	Frequency
User A	0.1534	25.68	3354
User B	0.3369	6.07	3369
User C	0.7034	5.94	3557
User D	0.7001	16.08	9725

Acknowledgments. — We thank Ming Tang for the valuable discussion. This work is jointly supported by the NNSFC(Grant Nos.90924011, 60933005, 61004102, 11105025), China Postdoctoral Science Foundation (Grant No. 20110491705), the Specialized Research Fund for the Doctoral Program of Higher Education (Grant No. 20110185120021). ZDZ appreciates the financial support of the Fundamental Research Funds for the Central Universities (Grant No. ZYGX2011YB024).

REFERENCES

- [1] BARABÁSI A.-L., *IEEE Contr. Syst. Mag.*, **27** (2007) 33.
- [2] ZHOU T., HAN X.-P. and WANG B.-H., *Towards the understanding of human dynamics* (Singapore: World Scientific Publishing) 2008.
- [3] BARABÁSI A.-L., *Nature*, **435** (2005) 207.
- [4] OLIVEIRA J. G. and BARABÁSI A.-L., *Nature*, **437** (2005) 1.
- [5] CANDIA J., GONZÁLEZ M. C., WANG P., SCHOENHARL T., MADEY G. and BARABÁSI A.-L., *Physica A*, **41** (2008) 224015.
- [6] DEZSÖ Z., ALMAAS E., LUKÁCS A., RÁCZ B., SZAKADÁT I. and BARABÁSI A.-L., *Phys. Rev. E*, **73** (2006) 066132.
- [7] ZHOU T., KIET H. A. T., KIM B. J., WANG B.-H. and HOLME P., *Europhys. Lett.*, **82** (2008) 28002.
- [8] GONÇALVES B. and RAMASCO J., *Phys. Rev. E*, **78** (2008) 026123.
- [9] RADICCHI F., *Phys. Rev. E*, **80** (2009) 026118.
- [10] ZHAO Z.-D. and ZHOU T., *Physica A*, **391** (2012) 3308.
- [11] HAN X.-P., ZHOU T. and WANG B.-H., *New J. Phys.*, **10** (2008) 073010.
- [12] VÁZQUEZ A., *Physica A*, **373** (2007) 747.
- [13] OLIVEIRA J. and VÁZQUEZ A., *Physica A*, **388** (2009) 187.
- [14] MIN B., GOH K.-I. and KIM I.-M., *Phys. Rev. E*, **79** (2009) 056110.
- [15] WU Y., ZHOU C., XIAO J., KURTHS J. and SCHELLNHUBER H., *Proc. Natl. Acad. Sci. U. S. A.*, **107** (2010) 18803.
- [16] PENG C.-K., BULDYREV S., HAVLIN S., SIMONS M., STANLEY H. E. and GOLDBERGER A. L., *Phys. Rev. E*, **49** (1994) 1685.
- [17] PENG C.-K., MIETUS J., HAUSDORFF J. M., HAVLIN S., STANLEY H. E. and GOLDBERGER A. L., *Phys. Rev. Lett.*, **70** (1993) 1343.
- [18] PENG C.-K., HAVLIN S., STANLEY H. E. and GOLDBERGER A. L., *Chaos*, **5** (1995) 82.

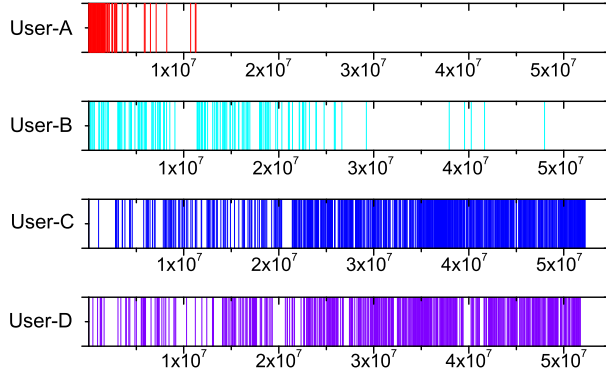


Fig. 5: (Color online) The four online active patterns of users corresponding to Tab. 2. The horizontal axis denotes time, and each vertical line indicates an individual event.

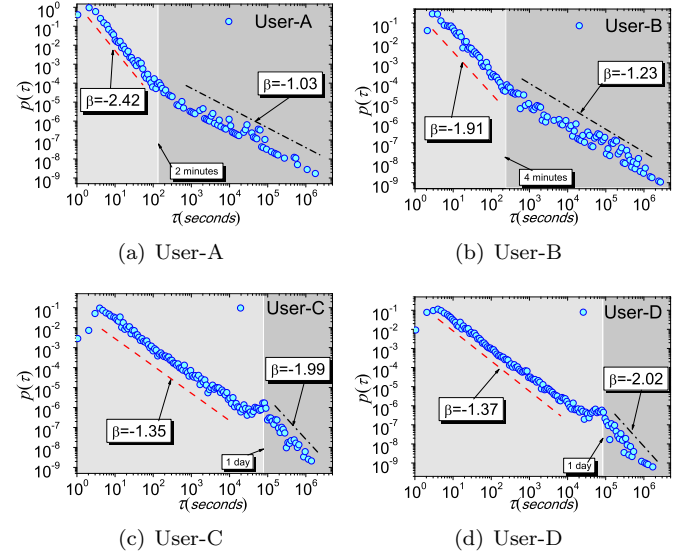


Fig. 6: Interevent time distributions.

- [19] IVANOV P. C., BUNDE A., AMARAL L. A. N., HAVLIN S., FRITSCH-YELLE J., BAEVSKY R. M., STANLEY H. E. and GOLDBERGER A. L., *Europhys. Lett.*, **48** (1999) 594.
- [20] IVANOV P. C., AMARAL L. A. N., GOLDBERGER A. L., HAVLIN S., ROSENBLUM M. G., STANLEY H. E. and STRUZIK, Z. R., *Chaos*, **11** (2001) 641.
- [21] ASHKENAZY Y. IVANOV P. C., HAVLIN S., PENG C.-K., GOLDBERGER A. L. and STANLEY H. E., *Phys. Rev. Lett.*, **86** (2001) 1900.
- [22] LIU Y., GOPIKRISHNAN P., CIZEAU P., MEYER M., PENG C.-K. and STANLEY H. E., *Phys. Rev. E*, **60** (1999) 1390.
- [23] GOPIKRISHNAN P., PLEROU V., GABAIX, X and STANLEY H. E., *Phys. Rev. E*, **62** (2000) R4493.
- [24] YAMASAKI K., MUCHNIK L., HAVLIN S., BUNDE A and STANLEY H. E., *Proc. Natl. Acad. Sci. U. S. A.*, **102** (2005) 9424.
- [25] IVANOV P. C., YUEN A., PODOBNIK B. and LEE Y.-K., *Phys. Rev. E*, **69** (2004) 056107.
- [26] JIANG Z.-Q., CHEN W, and ZHOU W.-X., *Physica A*, **388** (2009) 433.
- [27] TANG M. and LIU Z. H., *Physica A*, **387** (2008) 1361.
- [28] CAI S.-M., FU Z.-Q. ZHOU T., GU J. and ZHOU P.-L., *Europhys. Lett.*, **87** (2009) 68001.
- [29] LEVITIN D., CHORDIA P. and MENON V., *Proc. Natl. Acad. Sci. U. S. A.*, **109** (2012) 3716.
- [30] RYBSKI D., BULDYREV S. V., HAVLIN S., LILJEROS F. and MAKSE H. A., *Proc. Natl. Acad. Sci. U. S. A.*, **106** (2009) 12640.
- [31] RYBSKI D., BULDYREV S. V., HAVLIN S., LILJEROS F. and MAKSE H. A., *Arxiv preprint arXiv:1002.0216*, (2010) .
- [32] RYBSKI D., BULDYREV S. V., HAVLIN S., LILJEROS F. and MAKSE H. A., *Arxiv preprint arXiv:1205.1628*, (2012) .
- [33] COSTA M., GOLDBERGER A. and PENG C.-K., *Phys. Rev. Lett.*, **89** (2002) 068102.
- [34] COSTA M., GOLDBERGER A. L. and PENG C.-K., *Phys. Rev. Lett.*, **95** (2005) 198102.
- [35] MCINTOSH A. R., KOVACEVIC N. and ITIER R. J., *PLoS Computat. Biol.*, **4** (2008) e1000106.
- [36] ESCUDERO J, ABÁSULO D., HORNERO R., ESPINO P., and LÓPEZ, M., *Physiol. Meas.*, **27** (2006) 1091.
- [37] LI Z, ZHANG Y. K., *Stoch. Environ. Res. Risk Assess.*, **22** (2008) 507.
- [38] GOH K.-I. and BARABÁSI A.-L., *Europhys. Lett.*, **81** (2008) 48002.
- [39] HU K., IVANOV P., CHEN Z., CARPENA P. and STANLEY H. E., *Phys. Rev. E*, **64** (2001) 011114.
- [40] CHEN Z., IVANOV P., HU K. and STANLEY H. E., *Phys. Rev. E*, **65** (2002) 041107.
- [41] RICHMAN J. S. and MOORMAN J. R., *Am. J. Physiol.*, **278** (2000) H2039.
- [42] GHOSHAL G. and HOLME P., *Physica A*, **364** (2006) 603.